

# Integration of a heterogeneous compute resource in the ATLAS workflow

*Felix Bühner*<sup>1</sup>, *Anton Game*<sup>1,2,\*</sup>, *Benoit Roland*<sup>1</sup>, *Ulrike Schnoor*<sup>1,3</sup>, and *Markus Schumacher*<sup>1</sup> on behalf of the ATLAS Collaboration

<sup>1</sup>Institute of Physics, University of Freiburg, Freiburg, Germany

<sup>2</sup>Computing Center, University of Freiburg, Freiburg, Germany

<sup>3</sup>now at CERN, Geneva, Switzerland

**Abstract.** With the ever-growing amount of data collected with the experiments at the Large Hadron Collider (LHC), the need for computing resources that can handle the analysis of this data is also rapidly increasing. This increase will even be amplified after upgrading to the High Luminosity LHC [1]. High-Performance Computing (HPC) and other cluster computing resources provided by universities can be useful supplements to the resources dedicated to the experiment as part of the Worldwide LHC Computing Grid (WLCG) for data analysis and production of simulated event samples. Freiburg is operating a combined Tier2/Tier3, the ATLAS-BFG [2]. The shared HPC cluster "NEMO" at the University of Freiburg has been made available to local ATLAS [3] users through the provisioning of virtual machines incorporating the ATLAS software environment analogously to the bare metal system of the local ATLAS Tier2/Tier3 centre. In addition to the provisioning of the virtual environment, the on-demand integration of these resources into the Tier3 scheduler in a dynamic way is described. In order to provide the external NEMO resources to the user in a transparent way, an intermediate layer connecting the two batch systems is put into place. This resource scheduler monitors requirements on the user-facing system and requests resources on the backend-system.

## 1 Introduction

The analysis of collision data collected at the LHC and simulation of events is primarily done at 2 Tier0, 13 Tier1 and 160 Tier2 sites within the WLCG [4]. Tier3 components at sites provide resources also for local groups. WLCG clusters are mostly set up for High Throughput Computing (HTC) to get as much compute power as possible. High Performance Computing (HPC) clusters, as provided by universities and other institutions, sometimes even co-located at the same sites, may be used for HTC-like workflows to extend the capacities of the existing WLCG resources. However, there are no standard interfaces or abstraction layers available to easily cross-link clusters with different HPC/HTC setups, different resource managers or different login schemes. Therefore, configuring clusters to distribute workloads to several schedulers and share monitoring information is not a straight-forward task.

---

\*e-mail: [anton.game@physik.uni-freiburg.de](mailto:anton.game@physik.uni-freiburg.de)

In order to achieve the on-demand scheduling of resources on one cluster due to requirements on a different cluster, an intermediate layer, or resource scheduler, is put in place.

From its primary concept, the NEMO HPC cluster was designed to provide a full virtualization solution base on OpenStack [5] that enables users to spawn virtual machines (VMs) with a pre-configured image - so-called virtual research environments (VREs) [6]. These VMs are requested by sending a wrapper-job to the NEMO batch-system (MOAB [7]), which is queued in the same way as other jobs by NEMO users. When the wrapper-job starts, a virtual machine is spawned on the OpenStack instance. The lifetime of the VM is defined by the walltime of the MOAB job. After start-up and some initial checks, the VM is incorporated in the front-end scheduler on the ATLAS-BFG as an additional resource. This resource is in turn used to run the jobs that triggered the start of the VM in the first place. All of these mechanisms are completely transparent to the user of the ATLAS-BFG.

In the following, we describe how these virtual resources are integrated into the ATLAS Tier2/Tier3 cluster. NEMO and the Tier2/Tier3 are utilizing different batch systems. On the user-facing (or frontend) side, SLURM [8] is used as a scheduler, while the NEMO cluster (backend) runs a combination of MOAB & Torque. We use ROCED [9], developed at the Karlsruhe Institute of Technology (KIT) to schedule resources. ROCED is used already to integrate NEMO resources into the HTCondor [10] system at KIT. Due to the modular architecture of ROCED, it can also be used for connecting the two systems in Freiburg. To do so, a new component monitoring the SLURM queue on the frontend has been developed.

The performance of the system is being measured using several different benchmarks. These benchmarks are also used to quantify the modification of the performance due to changes in the configuration of the resource scheduler and of the virtual machines being spawned. They will also be part of a future continuous monitoring effort in order to be able to detect changes in the submitted workloads. This monitoring and tuning effort will ensure a robust but also dynamic and efficient environment.

## 2 Challenges

The ATLAS research groups in Freiburg have very specific requirements to the operating system as well as the installed software. This is to ensure reliable scientific results across all grid sites of the WLCG.

Virtualization has been found to be a technology that can simplify the challenge to provide a specific environment on a range of different heterogeneous and changing platforms, especially in the context of particle physics [11].

Being only one of multiple user groups on a shared HPC system, especially the choice of operating system has to take into account considerations from all user groups as well as from the party operating the cluster. A fully virtualized environment, independent of the choices made on the HPC cluster itself, will give the best possible scope to implement a system, that looks and behaves in the same way as the non-virtualized Tier2/Tier3 cluster. This consistency between the two systems would also make it possible in the future to redirect ATLAS grid jobs submitted remotely to either NEMO or any other opportunistic resource as long as the resource provides the needed infrastructure to run the VM images. The VM images which are made available to OpenStack on NEMO have to be created and updated easily in an automatic procedure and have to fulfil the following requirements:

- Scientific Linux 6 [12] - current OS on the Tier2/Tier3 cluster

- Access to ATLAS software via the CERN virtual file system CVMFS [13]
- User environment from Tier3
- Access to both grid-aware datasets on the distributed storage system dCache [14] and the local NEMO parallel filesystem BeeGFS [15].

Since the VMs are completely self-contained, all features needed to monitor and benchmark the machine are independent of the two schedulers that are involved and can either be implemented on the VM itself or offloaded to the resource scheduler. In the future, this information will also be used for continuous monitoring of the robustness and performance of the system.

### 3 Installation

#### 3.1 Generation of the virtual machines

The VM template is generated with Packer [16] and is based on a Scientific Linux 6 netinstall image. The customization and configuration of the template is done with puppet [17] which is also used for the contextualisation of the non-virtualized worker nodes on the Tier2/Tier3 cluster. Changes in configuration are automatically picked up by both systems. The output of this procedure is a static image that can be uploaded to the OpenStack server and is directly available.

In order to simplify software configurations across grid sites, most commonly-used software packages for the HEP-workflow are distributed using CVMFS. The software distributed on CVMFS is managed centrally by the experiments, hosted on web servers and transferred to the worker nodes on demand.

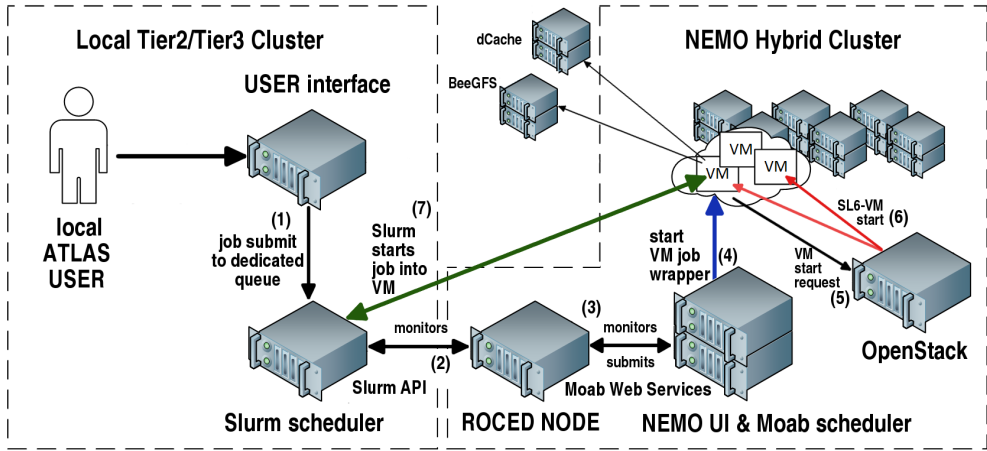
The VMs do not contain a predefined CVMFS cache and do not utilize the hard disk as persistent cache. A 256 MB RAM disk is being filled on demand from the local frontier squid [18] proxies. This circumvents an on-disk CVMFS installation on the host. This is a different approach to other solutions relying on access to software from CVMFS and using an on-disk cache like CernVM [19].

#### 3.2 Connection of front and backend batch systems

The biggest challenge for a smooth operation is the interconnection between the two different batch systems – SLURM on the Tier2/Tier3 frontend and MOAB on the backend (NEMO) through a resource scheduler.

The workflow is as follows: the resource scheduler monitors the frontend scheduler to which users send their workloads. The requirements are then compared to a list of available configurations on the NEMO HPC and the resource scheduler decides on the number of virtual machines for each configuration needed to fulfil the requirements. The appropriate number of batch jobs are sent to the backend scheduler, each spawning a virtual machine of the chosen type. The jobs that are used to start VMs in the OpenStack environment are regular user jobs on NEMO, which are started according to the availability of resources and the fair share of the NEMO user used to reproduce the fair share of the project. After startup of the VMs, they are integrated as additional resources into the SLURM scheduler and can then be used to process the user jobs queued in the frontend scheduler.

Figure 1 shows the general mode of operation. For now, SLURM in the Tier2/Tier3 cluster is set-up with separate partitions that reflect the user's affiliation to one of the local



**Figure 1.** Schematic view of how user jobs submitted to the frontend scheduler SLURM trigger jobs to start VMs on the OpenStack instance at NEMO.

ATLAS working groups. Each working group is in turn represented by a single user on NEMO, that is used to queue the jobs starting the VMs. By this mechanism, the fair share for the different areas of research using NEMO are incorporated into the workflow.

ROCED monitors these partitions. After user's job submission ROCED starts a wrapper job on the backend batch system, mapped to a single dedicated user per group. The wrapper job then triggers the start of a VM with a given configuration on OpenStack. As long as resources are requested and available on NEMO, additional virtual machines can be started. This mechanism leads to a dynamic extension of the amount of compute nodes and job slots available for physics analyses on the frontend system.

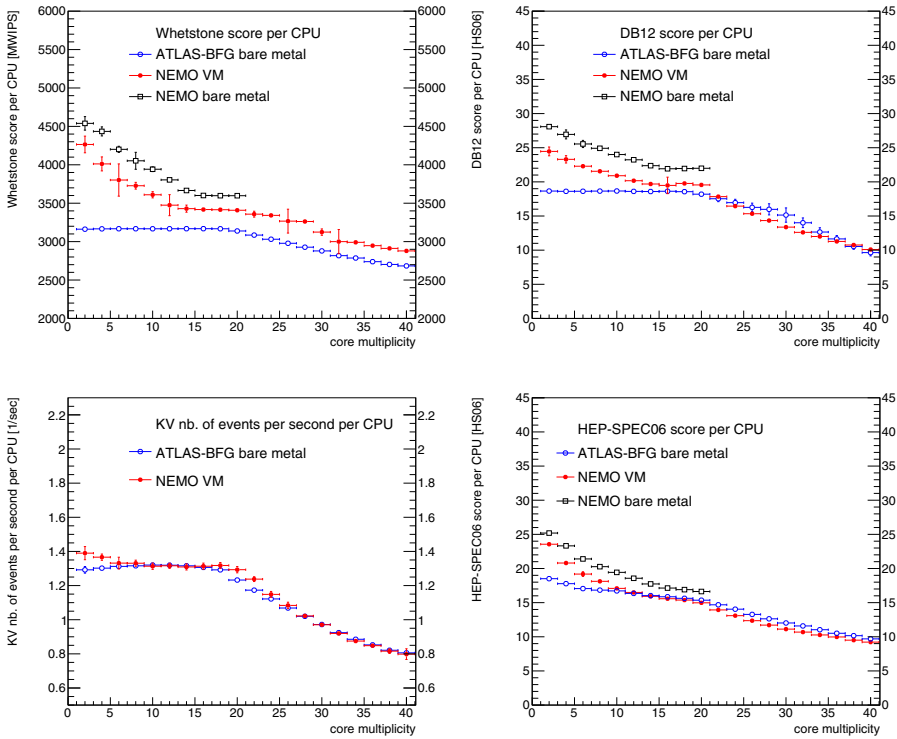
Before VMs are integrated into SLURM, a diagnosis-check is done to see whether all needed resources are available. After a successful check the VM is set to online in SLURM and jobs can be allocated.

## 4 Benchmarks

To understand the performance losses introduced by going to a fully virtualized environment, different benchmarks have been run. All benchmarks are carried out on the same hardware and the results obtained on the virtualized research environment are compared to the results running directly on hardware "bare metal") on both the Tier2/Tier3 and the NEMO cluster as well, to also assess the impact of different operating systems on the benchmark results.

In addition to the legacy HEP-SPEC06 (HS06) benchmark [20], the evaluation of the performance of the compute resources makes use of three benchmarking programs available in the CERN benchmark suite [21]:

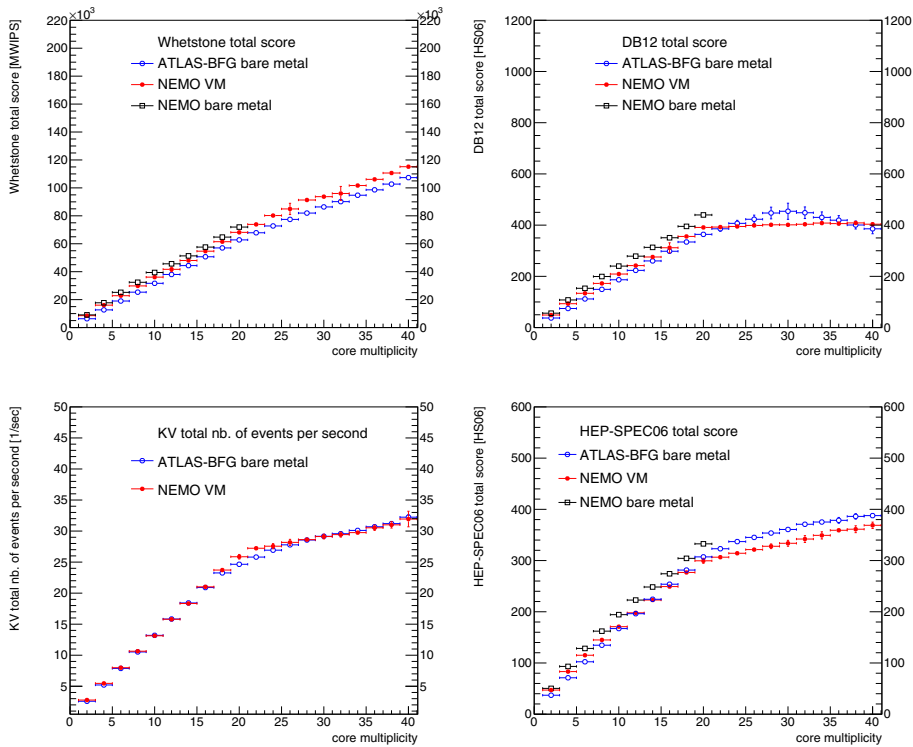
1. Dirac Benchmark 2012 DB12 [22]
2. Whetstone benchmark [23]
3. Kit Validation benchmark KV [24].



**Figure 2.** Score per CPU as a function of the core multiplicity for the Whetstone (top left), DB12 (top right), KV (bottom left) and HEP-SPEC06 (bottom right) benchmarks for the Tier2/Tier3 (ATLAS-BFG) running bare metal (blue open circles), the NEMO VMs (red full circles) and the NEMO bare metal (black open squares). The data points represent the average values of the benchmarks for each core multiplicity, and the vertical bars show the associated root-mean-squares. Horizontal error bars are only drawn for visibility and do not represent an uncertainty.

The DB12 and Whetstone programs are evaluating the performance of CPUs through floating-point arithmetic operations. One of the main differences between the two benchmarks resides in the variables used as input to the arithmetic operations: DB12 uses random numbers generated according to a Gaussian distribution, while Whetstone utilizes variables with predefined values. The KV benchmark runs the ATLAS software Athena[25] to simulate and reconstruct the interactions of muons in the detector of the ATLAS experiment.

As our primary target is to measure performance of CPUs in the context of High Energy Physics (HEP) applications, the KV benchmark constitutes a realistic payload, more suited to our goal than the DB12 and Whetstone software. The DB12 benchmark is measured in units of HS06 and therefore can be compared directly to the results from the HEP-SPEC06 benchmark. The Whetstone scores are expressed in Million of Whetstone Instructions Per Second (MWIPS), and the KV output provides the number of events produced per second. The different benchmarks are used to evaluate the performance of identical 20 cores Intel Xeon E5-2630 CPUs on the two different clusters: the Tier2/Tier3 cluster (ATLAS-BFG)



**Figure 3.** Total score as a function of the core multiplicity for the Whetstone (top left), DB12 (top right), KV (bottom left) and HEP-SPEC06 (bottom right) benchmarks for the Tier2/Tier3 (ATLAS-BFG) running bare metal (blue open circles), the NEMO VMs (red full circles) and the NEMO bare metal (black open squares). The data points represent the average values of the benchmarks for each core multiplicity, and the vertical bars show the associated root-mean-squares. Horizontal error bars are only drawn for visibility and do not represent an uncertainty.

and the shared HPC cluster (NEMO). The performance has been evaluated on three different configurations: the Tier2/Tier3 and NEMO HPC clusters running both on bare metal and the virtual machines running on the NEMO HPC cluster – except for the KV benchmark, which currently cannot be run on NEMO bare metal nodes.

On the Tier2/Tier3, hyperthreading (HT) technology is activated and the number of cores that can be used is higher by a factor of two with respect to the physical number of CPU cores available. For the virtual machines, an arbitrary number of CPU cores can be requested. The operating system used is Scientific Linux 6 in both cases. The NEMO bare metal has no HT activated due to the more general use case of the system, and uses CentOS7 [26] as operating system. The scores of the HEP-SPEC06, DB12, Whetstone and KV benchmarks have been determined for these three configurations as a function of the number of cores actually used by the benchmarking processes. This number ranges from 2 to 40 for the Tier2/Tier3 bare metal and for the VMs running on the NEMO cluster, for which HT is enabled, and from 2 to 20 for the NEMO bare metal, for which HT is not implemented. The results have been determined by step of two core units. The benchmarks have been run 20 times for each core multiplicity value, and the means and root-mean-squares (RMS) of the corresponding

distributions have been extracted.

The scores per CPU and the total scores are presented in Figures 2 and 3 respectively, for the four benchmarks and the three configurations considered, except for the KV software for which the NEMO bare metal results are not yet available. This latter benchmark retrieves the ATLAS software from the CVMFS file system, which is presently not available for the NEMO bare metal.

A decrease of the Whetstone, DB12 and HEP-SPEC06 scores per CPU is observed in Figure 2 for increasing values of the core multiplicity for the NEMO VMs and NEMO bare metal configurations. This observed decrease is significantly less pronounced for the Tier2/Tier3 bare metal results. There, the scores per CPU remain constant until the maximum number of physical cores is reached for all but the HS06 benchmark. The behaviour for the hyperthreaded region between 20 and 40 cores is similar for the two configurations. The KV results exhibit a similar behaviour for both the Tier2/Tier3 bare metal and the NEMO VMs, with a constant number of events produced per second per CPU below the maximum number of physical cores and a decrease of the performance afterwards. This behaviour is the closest to the pattern expected for an ideal CPU benchmark: a constant CPU performance per physical core and a decrease in the region where HT is active. The Whetstone score per CPU at a core multiplicity of 20, the maximum number of physical cores available, is considered as an illustrative example of the benchmark behaviours on the three different configurations. An increase of the CPU performance by the order of 5% is observed when going from the Tier2/Tier3 bare metal to the NEMO VMs, while going from the NEMO VMs to the NEMO bare metal leads to a further increase of performance of the order of 5% as well.

A continuously increasing total score is observed in Figure 3 for the Whetstone benchmark on the three different configurations, while the DB12, KV and HEP-SPEC06 results are characterized by a flattening increase or a constant behaviour once the maximum number of physical cores has been reached. The Whetstone benchmark provides higher CPU performance in the HT region in comparison to the scores obtained with the three other benchmarks. The scores obtained with the KV and HEP-SPEC06 benchmarks indicate an increase of the CPU performance by 15 to 20% when going from the maximum number of physical cores to the upper edge of the HT region, while the Whetstone scores exhibit a larger increase of the order of 60%. The Tier2/Tier3 bare metal and the VMs running on the NEMO cluster share the same configuration in terms of hardware, operating system and hyper-threading. A given benchmark should therefore exhibit a similar behaviour for both configurations. The KV benchmark, besides being the more realistic estimator of the CPU performance in the context of High Energy Physics applications, is the only benchmark for which this expectation is observed. The behaviours of the different benchmarks still need to be studied in more details, in order to fully understand the impact of the operating system, hyper-threading and virtualization on the CPU performance.

## 5 Summary

The HPC cluster NEMO has successfully been integrated into the workflow of local users in Freiburg running ATLAS data analysis jobs. This has been achieved in a transparent way using full virtualization on NEMO and utilizing the resource scheduler ROCED for the integration of the virtual resources into the frontend scheduler. The system is in production since fall 2017.

First performance tests using different benchmarks show some degradation in performance of the virtual machines compared to running on bare metal. The differences are within the

expected ranges when using different operating systems and are significantly reduced when going from pure CPU benchmarks like Whetstone or DB12 to benchmarks more closely related to high energy particle physics analysis like HS06 or KV.

The continuous benchmarking effort will ensure a stable and efficient environment. Integrating the results of the benchmarks into the resource scheduling process can enable cluster administrators to test different configurations, leading to a more efficient usage of the provided resources.

## 6 Acknowledgements

The research is supported by the Ministerium für Wissenschaft, Forschung und Kunst Baden-Württemberg in the project MWK 34-7547.221 “ViCE Virtual Open Science Collaboration Environment” and by the Bundesministerium für Bildung und Forschung in the project 05H15VFCA1 “Higgs-Physik mit dem und Grid-Computing für das ATLAS-Experiment am LHC”.

## References

- [1] G. Apollinari, O. Bruening, T. Nakamoto, L. Rossi (2017), [arXiv:1705.08830](https://arxiv.org/abs/1705.08830)
- [2] R. Backofen, H.G. Borrmann, W. Deck, A. Dedner, L. De Raedt, K. Desch, M. Diesmann, M. Geier, A. Greiner, W. R. Hess et al., *Praxis der Informationsverarbeitung und Kommunikation* **29**, 81 (2006)
- [3] G. Aad et al. (ATLAS), *JINST* **3**, S08003 (2008)
- [4] J. Shiers, *Computer Physics Communications* **177**, 219 (2007), proceedings of the Conference on Computational Physics 2006
- [5] OpenStack Foundation, *OpenStack (Newton)* (2010), <https://www.openstack.org/>
- [6] D. von Suchodoletz, B. Wiebelt, K. Meier, M. Janczyk, *Flexible HPC: bwForCluster NEMO*, in *Proceedings of the 3rd bwHPC-Symposium: Heidelberg 2016* (heiBOOKS, 2017), <http://books.ub.uni-heidelberg.de/heibooks/reader/download/308/308-4-79237-1-10-20171002.pdf>
- [7] Adaptive Computing, *MOAB HPC SUITE* (2014), <http://www.adaptivecomputing.com/products/hpc-products/moab-hpc-suite-grid-option/>
- [8] M.A. Jette, A.B. Yoo, M. Grondona, *SLURM: Simple Linux Utility for Resource Management*, in *In Lecture Notes in Computer Science: Proceedings of Job Scheduling Strategies for Parallel Processing (JSSPP) 2003* (Springer-Verlag, 2002), pp. 44–60
- [9] G. Erli, F. Fischer, G. Fleig, M. Giffels, T. Hauth, G. Quast, M. Schnepf, J. Heese, K. Leppert, J.A. de Pedro et al., *Journal of Physics: Conference Series* **898**, 052021 (2017)
- [10] University of Wisconsin – Madison, *HTCondor* (2018), <https://research.cs.wisc.edu/htcondor/>
- [11] P. Buncic, C. Aguado Sánchez, J. Blomer, A. Harutyunyan, M. Mudrinic, *The European Physical Journal Plus* **126**, 13 (2011)
- [12] Fermilab and CERN, *Scientific Linux release 6.8 (Carbon)* (2011), <http://www.scientificlinux.org/>
- [13] J. Blomer, G. Ganis, N. Hardi, R. Popescu, *Delivering LHC Software to HPC Compute Elements with CernVM-FS*, in *High Performance Computing*, edited by J.M. Kunkel, R. Yokota, M. Tauffer, J. Shalf (Springer International Publishing, Cham, 2017), pp. 724–730, ISBN 978-3-319-67630-2



- [14] A.P. Millar, G. Behrmann, C. Bernardt, P. Fuhrmann, D. Litvintsev, T. Mkrtchyan, A. Petersen, A. Rossi, K. Schwank, *Journal of Physics: Conference Series* **513**, 042033 (2014)
- [15] The Fraunhofer-Institut für Techno- und Wirtschaftsmathematik ITWM, *BeeGFS* (2014), <https://www.beegfs.io>
- [16] HashiCorp, *Packer 1.2.4* (2013), <https://www.packer.io/>
- [17] Puppet, *puppet 3.8.7* (2005), <https://puppet.com>
- [18] B. Blumenfeld, D. Dykstra, L. Lueking, E. Wicklund, *Journal of Physics: Conference Series* **119**, 072007 (2008)
- [19] P. Buncic, C. Aguado-Sanchez, J. Blomer, A. Harutyunyan, *Journal of Physics: Conference Series* **331**, 052004 (2011)
- [20] HEPiX Benchmarking Working Group, *HEP SPEC06 (HS06) benchmark* (2006), <https://w3.hepox.org/benchmarking.html>
- [21] M. Alef, C. Cordeiro, A. De Salvo, A. Di Girolamo, L. Field, D. Giordano, M. Guerri, F.C. Schiavi, A. Wiebalck, *J. Phys. Conf. Ser.* **898**, 092056 (2017)
- [22] R. Graciani, A. McNab, *Dirac benchmark 2012* (2012), [gitlab.cern.ch/mcnab/dirac-benchmark/tree/master](https://gitlab.cern.ch/mcnab/dirac-benchmark/tree/master)
- [23] H. Curnow, B. Wichman, *Computer Journal* **19**, 43 (1976)
- [24] A.D. Salvo, F. Brasolin, *Journal of Physics: Conference Series* **219**, 042037 (2010)
- [25] P. Calafiura, W. Lavrijsen, C. Leggett, M. Marino, D. Quarrie, *The Athena control framework in production, new developments and lessons learned*, in *Computing in high energy physics and nuclear physics. Proceedings, Conference, CHEP'04, Interlaken, Switzerland, September 27-October 1, 2004* (2005), pp. 456–458
- [26] The CentOS Project, *CentOS Linux release 7.4.1708 (Core)* (2017), <https://www.centos.org/>